# TODD GARNER

DS 6306 WEEK 6 PART 1

FEBRUARY 87. 2023

# READ IN THE DATA AND CLEAN – PART 1, #1

- The site www.opendatasoft.com proved challenging so I found the .csv files on kaggle.com and also in the class GitHub repo.

```r
96  ```{r}
97  # Load the Titanic dataset
98  Titanic <- read.csv(file.choose(), header = TRUE)
99  library(dplyr)
100 Titanic_new <- Titanic %>% select(Survived, Pclass, Age)
101 Titanic_exp <- Titanic %>% select(Pclass, Age)
102 Titanic_new_final <- Titanic_new[complete.cases(Titanic_new), ]
103 Titanic_final <- Titanic_exp[complete.cases(Titanic_exp), ]
104 ```
```

- Two data sets:  training and test.  Removed the NA's and selected the three pertinent columns.

- Comment:  I used ChatGPT and found it distracted me more than it assisted me.  I may use it in a different manner than writing the code base.

# CREATE THE TRAINING/TESTING SETS, USE THEM AND CONFUSION MATRIX – PART 1, #2

```
114
115  set.seed(123)
116  train_new_index = sample(1:nrow(Titanic_new_final), size = nrow(Titanic_new_final)*0.8, replace = FALSE)
117  train_new = Titanic_new_final[train_new_index, ]
118  test_new = Titanic_new_final[-train_new_index, ]
119
120  set.seed(123)
121  train_index = sample(1:nrow(Titanic_final), size = nrow(Titanic_final)*0.8, replace = FALSE)
122  train = Titanic_final[train_index, ]
123  test = Titanic_final[-train_index, ]
124
125  library(class)
126  library(caret)
127  knn.23 <- knn(train, test, train_new$Survived, k = 23)
128  knn.24 <- knn(train, test, train_new$Survived, k = 24)
129  table(test_new$Survived, knn.23)
130  confusionMatrix(table(test_new$Survived, knn.23))
131
132  table(test_new$Survived, knn.24)
133  confusionMatrix(table(test_new$Survived, knn.24))
134
```

- Data set: training -> 571          testing -> 143  (177 NA's removed)

- Transforming in KNN and Confusion Matrix (results on next slide)

# CONFUSION MATRIX AND STATISTICS

```
Confusion Matrix and Statistics

    knn.24
     0  1
  0 69 10
  1 46 18

                Accuracy : 0.6084
                  95% CI : (0.5233, 0.6889)
     No Information Rate : 0.8042
     P-Value [Acc > NIR] : 1

                   Kappa : 0.1634

 Mcnemar's Test P-Value : 2.91e-06

             Sensitivity : 0.6000
             Specificity : 0.6429
          Pos Pred Value : 0.8734
          Neg Pred Value : 0.2813
              Prevalence : 0.8042
          Detection Rate : 0.4825
    Detection Prevalence : 0.5524
       Balanced Accuracy : 0.6214

        'Positive' Class : 0
```

Having seen the answers, I know now that my solution was off the mark. I'm impressed with your code. Specific and on target. I have much to learn.

# AGE AND PROBABILITY OF SURVIVAL, PART 1, #3

```r
142
143 ```{r}
144 knn.23 <- knn(train, test, train_new$Survived, k = 23)
145 knn.24 <- knn(train, test, train_new$Survived, k = 24)
146
147 table(test_new$Survived, knn.23)
148 confusionMatrix(table(test_new$Survived, knn.23))
149 table(test_new$Survived, knn.24)
150 confusionMatrix(table(test_new$Survived, knn.24))
151 test_new <- data.frame(Pclass = 1, Age = 59)
152 test_new.2 <- data.frame(Pclass = 2, Age = 59)
153 test_new.3 <- data.frame(Pclass = 3, Age = 59)
154 knn.todd_1 <- knn(train, test_new, train_new$Survived, k = 24)
155 knn.todd_2 <- knn(train, test_new.2, train_new$Survived, k = 24)
156 knn.todd_3 <- knn(train, test_new.3, train_new$Survived, k = 24)
157
```

```
[1] 0
Levels: 0 1
[1] 0
Levels: 0 1
[1] 0
Levels: 0 1
```

Through my research, I found one suggestion on the value to select for k. It was suggested that k can be derived from the square root of the size of the data set.

The result is between k = 23 and k = 24. I've employed them both

I created 3 data sets with my Age (59) in it along with the Pclass. 1st, 2nd, 3rd.

Given the data, it looks bad for someone my age, regardless of the class.

# USE OF THE TRAINING SET AND THE TEST SET – PART 1 #4

- Utilizing the prior training set (571 observations), I've loaded in the test_set.csv. This is where I got stuck and moved to Part 2.

```r
165 ```{r}
166 Titanic_418 <- read.csv(file.choose(), header = TRUE)
167 Titanic_418_new <- Titanic_418 %>% select(Pclass, Age)
168 Titanic_418_final <- Titanic_418_new[complete.cases(Titanic_418_new), ]
169 knn.418_final <- knn(train, Titanic_418_final, train_new$Survived, k = 18)
170 train_new$Survived <- as.factor(train_new$Survived)
171 ```
```

# PART 2 # A. IRIS DATA SET – 70/30 TRAINING SET, PLOT XAXIS

```r
18  ```{r}
19  library(tidyverse)
20  iris
21  View(iris)
22  splitPerc = .70
23  train_init = sample(1:dim(iris)[1],round(splitPerc * dim(iris)[1]))
24  train = iris[train_init,]
25  test = iris[-train_init,]
26  ```
```

Description: df [150 x 5]

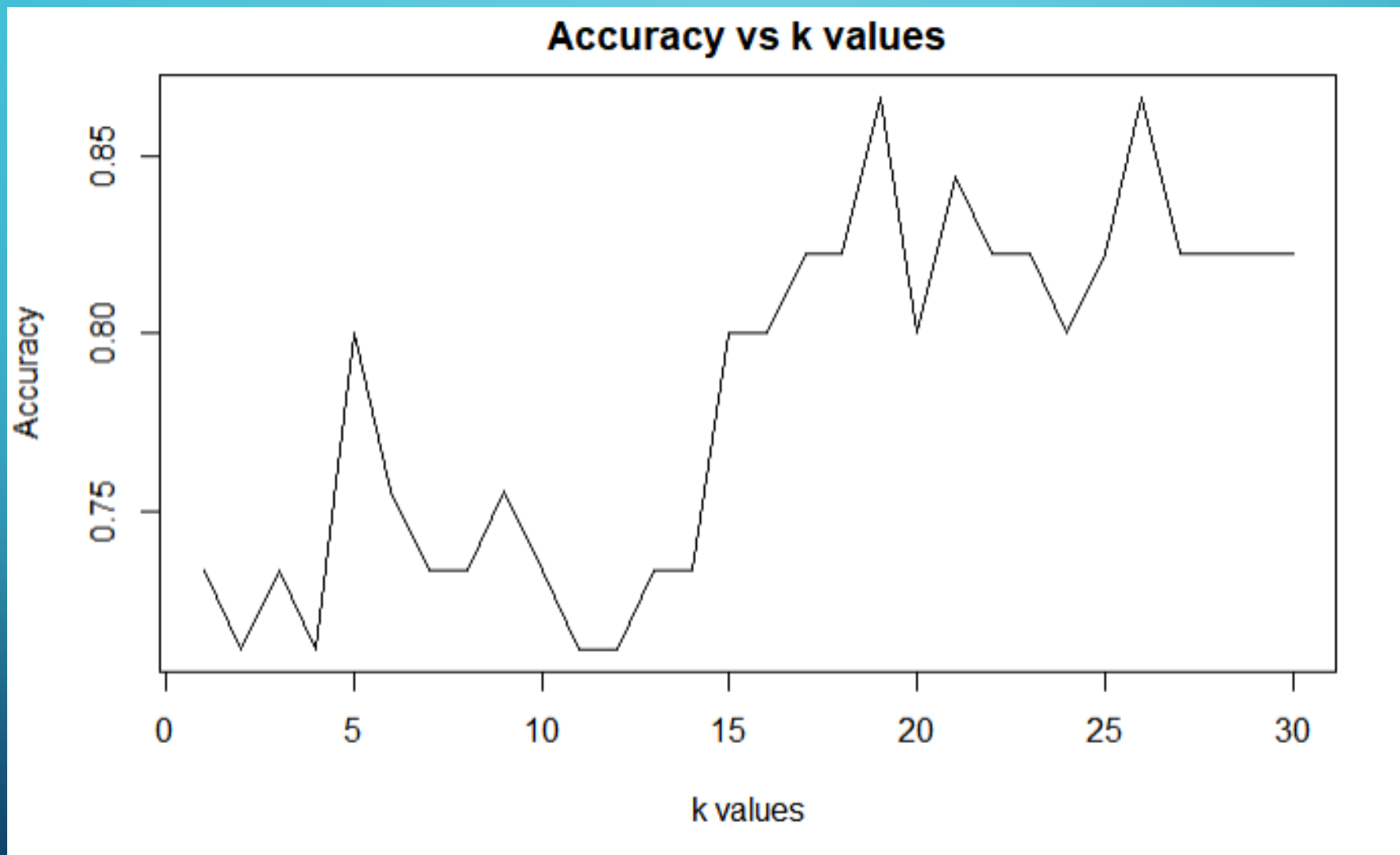| Sepal.Length <dbl> | Sepal.Width <dbl> | Petal.Length <dbl> | Petal.Width <dbl> | Species <fctr> |
|---|---|---|---|---|
| 5.1 | 3.5 | 1.4 | 0.2 | setosa |
| 4.9 | 3.0 | 1.4 | 0.2 | setosa |
| 4.7 | 3.2 | 1.3 | 0.2 | setosa |
| 4.6 | 3.1 | 1.5 | 0.2 | setosa |
| 5.0 | 3.6 | 1.4 | 0.2 | setosa |
| 5.4 | 3.9 | 1.7 | 0.4 | setosa |
| 4.6 | 3.4 | 1.4 | 0.3 | setosa |
| 5.0 | 3.4 | 1.5 | 0.2 | setosa |
| 4.4 | 2.9 | 1.4 | 0.2 | setosa |
| 4.9 | 3.1 | 1.5 | 0.1 | setosa |

```r
32  iris
33  library(class)
34  library(FNN)
35  library(caret)
36  library(ISLR)
37  set.seed(6)
38  splitPerc = .70
39  trainIndices = sample(1:dim(iris)[1],round(splitPerc * dim(iris)[1]))
40  train = iris[trainIndices,]
41  test = iris[-trainIndices,]
42  iris %>% ggplot(aes(x = Sepal.Length,Sepal.Width,color = Species)) + geom_point()
43  View(train$Species)
44  # k = 5
45  classifications = knn(train[,c(1,2)],test[,c(1,2)],train$Species, prob = TRUE, k = 5)
46  table(test$Species,classifications)
47  confusionMatrix(table(test$Species,classifications))
```

Iris – Split the data into 70/30 train/test sets, plot the results (plots on the following slide), classify and employ the confusion matrix.

```r
49  accs = data.frame(accuracy = numeric(30), k = numeric(30))
50  nrow(train)
51  for(i in 1:30)
52  {
53  classifications = knn(train[,c(1,2)],test[,c(1,2)],train$Species, prob = TRUE, k = i)
54  table(test$Species,classifications)
55  CM = confusionMatrix(table(test$Species,classifications))
56  accs$accuracy[i] = CM$overall[1]
57  accs$k[i] = i
58  }
59  plot(accs$k,accs$accuracy, type = "l")
```

Create a "for loop" to test the k values to find the best possible fit.

# DETERMINATION OF BEST K VALUE | ANSWER = 18



Accuracy vs k values

# CONFUSION MATRIX RESULTS

```
    setosa          13           0           0
    versicolor       1          11           4
    virginica        0           3          13

Overall Statistics

               Accuracy : 0.8222
                 95% CI : (0.6795, 0.92)
    No Information Rate : 0.3778
    P-Value [Acc > NIR] : 1.262e-09

                  Kappa : 0.7327

 Mcnemar's Test P-Value : NA

Statistics by Class:

                     Class: setosa Class: versicolor Class: virginica
Sensitivity                 0.9286            0.7857           0.7647
Specificity                 1.0000            0.8387           0.8929
Pos Pred Value              1.0000            0.6875           0.8125
Neg Pred Value              0.9688            0.8966           0.8621
Prevalence                  0.3111            0.3111           0.3778
Detection Rate              0.2889            0.2444           0.2889
Detection Prevalence        0.2889            0.3556           0.3556
Balanced Accuracy           0.9643            0.8122           0.8288
[1] 105
```

# BONUS QUESTION: REPEAT PRIOR ANALYSIS WITH LOOCV METHODOLOGY - RESULTS

```
        setosa          50          0          0
        versicolor       0         31         12
        virginica        0         19         38

Overall Statistics

              Accuracy : 0.7933
                95% CI : (0.7197, 0.8551)
    No Information Rate : 0.3333
    P-Value [Acc > NIR] : < 2.2e-16

                 Kappa : 0.69

 Mcnemar's Test P-Value : NA

Statistics by Class:

                     Class: setosa Class: versicolor Class: virginica
Sensitivity                 1.0000            0.6200           0.7600
Specificity                 1.0000            0.8800           0.8100
Pos Pred Value              1.0000            0.7209           0.6667
Neg Pred Value              1.0000            0.8224           0.8710
Prevalence                  0.3333            0.3333           0.3333
Detection Rate              0.3333            0.2067           0.2533
Detection Prevalence        0.3333            0.2867           0.3800
Balanced Accuracy           1.0000            0.7500           0.7850
[1] 105
```

These results are lower slightly than the k = 18 confusion matrix results on the prior page.

I think this highlights that LOOCV is great to use on small-ish data sets. On data sets with hundreds of values/observations, it appears that LOOCV is not superior to determining an external (XX%/XX%) versus an internal (LOOCV) methodology.